



Contaminating viral sequences in high-throughput sequencing viromics a linkage study of 700 sequencing libraries

Asplund, Maria; Kjartansdóttir, Kristín Rós; Mollerup, Sarah; Vinner, Lasse; Fridholm, Helena; Herrera, José A R; Friis-Nielsen, Jens; Hansen, Thomas Arn; Jensen, Randi Holm; Nielsen, Ida Broman; Richter, Stine Raith; Rey-Iglesia, Alba; Matey-Hernandez, Maria Luisa; Alquezar-Planas, David E; Olsen, Pernille V S; Sicheritz-Pontén, Thomas; Willerslev, Eske; Lund, Ole; Brunak, Søren; Mourier, Tobias; Nielsen, Lars Peter; Izarzugaza, Jose M G; Hansen, Anders Johannes

Published in:
Clinical Microbiology and Infection

DOI:
[10.1016/j.cmi.2019.04.028](https://doi.org/10.1016/j.cmi.2019.04.028)

Publication date:
2019

Document version
Publisher's PDF, also known as Version of record

Document license:
[CC BY-NC-ND](#)

Citation for published version (APA):
Asplund, M., Kjartansdóttir, K. R., Mollerup, S., Vinner, L., Fridholm, H., Herrera, J. A. R., Friis-Nielsen, J., Hansen, T. A., Jensen, R. H., Nielsen, I. B., Richter, S. R., Rey-Iglesia, A., Matey-Hernandez, M. L., Alquezar-Planas, D. E., Olsen, P. V. S., Sicheritz-Pontén, T., Willerslev, E., Lund, O., Brunak, S., ... Hansen, A. J. (2019). Contaminating viral sequences in high-throughput sequencing viromics: a linkage study of 700 sequencing libraries. *Clinical Microbiology and Infection*, 25(10), 1277-1285. <https://doi.org/10.1016/j.cmi.2019.04.028>



Original article

Contaminating viral sequences in high-throughput sequencing viromics: a linkage study of 700 sequencing libraries

M. Asplund^{1,*}, K.R. Kjartansdóttir¹, S. Møllerup¹, L. Vinner¹, H. Fridholm^{1,2}, J.A.R. Herrera^{3,4}, J. Friis-Nielsen⁴, T.A. Hansen¹, R.H. Jensen¹, I.B. Nielsen¹, S.R. Richter¹, A. Rey-Iglesia¹, M.L. Matey-Hernandez⁴, D.E. Alquezar-Planas¹, P.V.S. Olsen¹, T. Sicheritz-Pontén^{1,5}, E. Willerslev¹, O. Lund⁴, S. Brunak^{3,4}, T. Mourier¹, L.P. Nielsen², J.M.G. Izarzugaza⁴, A.J. Hansen^{1,**}

¹ Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Copenhagen, Denmark

² Department of Autoimmunology and Biomarkers, Statens Serum Institut, Copenhagen, Denmark

³ Disease Systems Biology Programme, Panum Institut, Copenhagen, Denmark

⁴ Department of Bio and Health Informatics, Technical University of Denmark, Lyngby, Denmark

⁵ Centre of Excellence for Omics-Driven Computational Biodiscovery, AIMST University, Kedah, Malaysia

ARTICLE INFO

Article history:

Received 3 January 2019

Received in revised form

12 April 2019

Accepted 18 April 2019

Available online 4 May 2019

Editor: L. Kaiser

Keywords:

Cluster

Contamination

High-throughput sequencing

Laboratory component

Metagenomic

Next-generation sequencing

Nucleic acid

Virome

Virus

ABSTRACT

Objectives: Sample preparation for high-throughput sequencing (HTS) includes treatment with various laboratory components, potentially carrying viral nucleic acids, the extent of which has not been thoroughly investigated. Our aim was to systematically examine a diverse repertoire of laboratory components used to prepare samples for HTS in order to identify contaminating viral sequences.

Methods: A total of 322 samples of mainly human origin were analysed using eight protocols, applying a wide variety of laboratory components. Several samples (60% of human specimens) were processed using different protocols. In total, 712 sequencing libraries were investigated for viral sequence contamination. **Results:** Among sequences showing similarity to viruses, 493 were significantly associated with the use of laboratory components. Each of these viral sequences had sporadic appearance, only being identified in a subset of the samples treated with the linked laboratory component, and some were not identified in the non-template control samples. Remarkably, more than 65% of all viral sequences identified were within viral clusters linked to the use of laboratory components.

Conclusions: We show that high prevalence of contaminating viral sequences can be expected in HTS-based virome data and provide an extensive list of novel contaminating viral sequences that can be used for evaluation of viral findings in future virome and metagenome studies. Moreover, we show that detection can be problematic due to stochastic appearance and limited non-template controls. Although the exact origin of these viral sequences requires further research, our results support laboratory-component-linked viral sequence contamination of both biological and synthetic origin. **M. Asplund, Clin Microbiol Infect 2019;25:1277**

© 2019 The Author(s). Published by Elsevier Ltd on behalf of European Society of Clinical Microbiology and Infectious Diseases. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Introduction

High-throughput sequencing (HTS) is an indispensable tool in life science research and clinical diagnostics [1,2] and facilitates the generation of massive amounts of DNA sequence information at acceptable costs within a short timeframe. The field of viromics has benefited from the rapid improvement of HTS technologies, as evidenced by major discoveries of novel viruses [3–9], some of

* Corresponding author. M. Asplund, Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Copenhagen, Denmark.

** Corresponding author. A.J. Hansen, Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Copenhagen, Denmark.

E-mail addresses: amasplund@snm.ku.dk (M. Asplund), ajhansen@snm.ku.dk (A.J. Hansen).

which have proven to be the cause of recent human epidemics [10]. Due to high sequence diversity, it has been challenging to identify novel viral genomes in clinical specimens using sequence-specific molecular methods such as PCR. HTS technologies provide an attractive alternative approach for virus discovery that require no previous knowledge about viral genomes. However, discovery of viruses using HTS also poses a number of challenges that must be accounted for in the interpretation of data. Sample preparation for HTS includes treatment with various laboratory components, also used for sample preparation in other non-HTS methods. Laboratory components have previously been documented to carry viral nucleic acid contamination [11–17]. Great caution is therefore essential when claiming disease association with a particular microorganism, to avoid incorrect conclusions, as in some unfortunate recent examples [18–22]. A better understanding of laboratory-component-derived contamination is needed.

Here, we systematically address the problem of nucleic acid contamination using HTS with focus on virus identification in clinical samples. We provide a comprehensive *in silico* characterization of contaminating viral sequences and their probable sources. More than 300 samples were analysed using eight overall methods applying an extensive variety of laboratory components. The original purpose of the investigation was to identify sample-derived viral sequences, findings described in Møllerup et al. (under review). In many cases (165 out of 274 human specimens), the same sample was processed using different laboratory protocols, resulting in several sequencing libraries per sample. Consequently, this study poses a unique opportunity for the characterization of common viral artefacts and contaminants in HTS metagenomic studies within clinical and other samples.

Methods

Ethics statement

The Regional Committee on Health Research Ethics and the National Committee on Health Research Ethics decided that ethical permission was not needed for collection and processing of these samples (case no. H-2-2012-FSP2 and 1304226) according to the Danish national legislation (Sundhedsloven). The samples used in this study were processed anonymously. All experiments were conducted according to the Declaration of Helsinki.

Samples

Samples consisted of 274 human specimens (32 different sample types, mostly of cancerous origin), 5 virus-spiked positive control samples and 43 non-template controls (NTCs) (see [Supplementary material, table s1](#)). Laboratory method development was not part of this study and positive controls were included to assess the bioinformatic pipeline.

Sample processing

To identify viral sequences within the samples, eight different overall methods were used; four DNA-focused methods and four RNA-focused methods (shotgun DNA and RNA, circular DNA enrichment, virion enrichment DNA and RNA, mRNA enrichment, retrovirus capture DNA and mRNA) (see [Supplementary material, Laboratory methods s1](#)). A total of 712 sequencing libraries were prepared and sequenced on the Illumina HiSeq 2000 platform with 2 × 100 bp paired-end sequencing. For sample processing, 54 laboratory reagents and utilities (laboratory components) were applied (see [Supplementary material, fig. s1 and table s2](#)). All samples were processed in the same laboratory.

Characterization of sequencing data

Paired-end sequencing reads were adapter trimmed and quality trimmed and merged. Reads mapping to the human reference genome (hg38), reads <30 nucleotides in length and low-complexity reads were excluded from further analysis. Remaining reads were assembled into larger contiguous sequences (contigs) from a combination of pairs, collapsed (merged overlapping pairs) and singleton reads. Default parameters were used for this purpose. Contigs and all human depleted and quality filtered reads were queried against the NCBI nucleotide database (nt) using BLASTn (MEGABLAST) [23] with a cut-off e-value of 10^{-3} . Contigs with no BLASTn hit were queried against the NCBI non-redundant protein database (nr) using BLASTx with the same cut-off e-value. For each characterized sequence the best hit was selected and taxonomically classified using the NCBI taxonomy database. All sequences with a viral classification were selected and sequences with the same viral taxID at the first level (species/strain) were clustered. Reads possibly occurring because of library misidentification, as a result of mixed sequencing clusters, referred to as bleedover [24], were considered. A bleedover ratio was calculated by dividing the viral read count of each viral sequence with the highest viral read count for the same viral sequence from different libraries sequenced on the same lane. Identified viral sequences with bleedover ratio <0.3% were removed. Cross-contamination from one sequencing run to another was not considered but could potentially also be present. Hosts of viruses were recovered from the NCBI taxonomy browser. Statistical analysis and visualization of data were performed using the software R v. 3.5.1 [25].

Association analysis

The identified viral sequences were correlated to laboratory components and sample types to detect possible sources of contamination. This was done using a positive one-tailed Fisher's exact test (significance level $\alpha = 0.05$ with Bonferroni correction).

Coverage analysis

A reference genome was selected for each viral sequence linked to a laboratory component (see [Supplementary material, table s3](#)). Using BOWTIE2 v. 2.2.5 [26] human depleted and quality filtered reads were mapped to viral reference genomes, applying global end-to-end and local alignment. Independently, the same reads were mapped to six manually selected algal chloroplast genomes. The alignments of reads to the reference genomes was visualized using CIRCOS (v0.67-7) [27], and an additional analysis of correlation to features based on mapping results was conducted. Cross-library genome coverage was investigated by summing the library-specific genome coverage of all libraries.

Results

A diversity of viral enrichment methods was applied to 279 samples of mainly human cancerous origin, resulting in 712 sequencing libraries (see [Table 1](#)).

A total of 56 728 213 824 sequencing reads were generated. After human depletion and quality filtering 2 953 972 594 reads and 1 381 107 contigs were characterized using BLAST. The results are summarized in [Table 2](#) (for library-specific information see [Supplementary material, table s1](#)).

Viral sequences linked to laboratory components

From BLAST of reads and contigs 2994 viral clusters were identified (see Methods). Of these, significant associations were

Table 1
Samples and libraries included in the study

Sample type	Samples	Shotgun DNA	Shotgun RNA	Virion enrichment		Circular DNA enrich.	Capture		mRNA enrich-ment	Libraries
				DNA	RNA		Retro-virus DNA	Retro-virus mRNA		
Basal cell carcinoma (cutaneous)	11	11		11	11	4	6			43
Mycosis fungoides (cutaneous)	11	11		11	11	10	11			54
Melanoma (cutaneous)	10	10		10	10	8				38
Oral cancer	10	12		10	10	10				42
Oral healthy	1					1				1
Vulvar cancer	3			3	4	3				10
Bladder cancer	7			8	9	5				22
Bladder cancer, urine	11		2			10				12
Colon cancer	32	12	11	3	3		6		6	41
Colon cancer, blood	8	8								8
Colon cancer, ascites	1	1					1			2
Colon healthy	2								2	2
Breast cancer (ductal)	10	10	10	9	13	8				50
Breast cancer (lobular)	10	10	9	10	10	7				46
Breast cancer, ascites	2	1	1	1	1	2				6
Testicular cancer (seminoma)	11	1		11	12					24
Testicular cancer (non-seminoma)	5	3		5	8					16
Testicular cancer (seminoma and non-seminoma)	4	1		4	4					9
AML	15		6	9	9	7				31
B-CLL	17		8	9	9	8				34
BCP-ALL	8			8	8	8				24
CML	20		10	10	10	10				40
T-ALL	20		9	11	11	9				40
Ovarian cancer, ascites	10	5	4	3	3	5				20
Pancreatic cancer, ascites	4	2	2				1			5
Optic neuritis, cerebrospinal fluid	4			4						4
Optic neuritis, plasma	4			4						4
Vasculitis	4			4						4
Gynaecological observation, ascites	1		1							1
Cell lines	18	12						6	6	24
Positive control	5	10					2			12
NTC				20	18	5				43
Total	279	120	73	178	174	120	27	6	14	712

Abbreviations: ALL, acute lymphoblastic leukaemia; AML, acute myeloid leukaemia; B-CLL, B-cell chronic lymphocytic leukaemia; BCP-ALL, B-cell precursor ALL; CML, chronic myeloid leukaemia; NTC, non-template control.

The table shows the number of samples for each sample type, the number of samples processed with the different laboratory methods, and the resulting number of libraries for each sample type (rightmost column) and laboratory method (bottom line).

found between 493 viral clusters and laboratory components, hereafter referred to as laboratory-component-associated (LCA) viral sequences (see Fig. 1 and Supplementary material, fig. s2a,b and table s4). Remarkably, 68% (62 521 069) of all viral reads were included in viral clusters linked to laboratory components (see Supplementary material, fig. s3a). For viral contigs this number was 74% (13 687 contigs). The majority of LCA viral sequences were non-human (see Supplementary material, fig. s3b), with 60% (296/493) being bacteriophages.

Some of the laboratory components showed high correlation when investigating the extent of simultaneous use (see Supplementary material, fig. s4), which can explain viral sequences showing significant association to multiple laboratory components. A particularly high proportion of viral sequences linked to laboratory components was seen for RNA-targeting overall methods (see Supplementary material, fig. s3c). Components used as part of RNA

methods (RNeasy MinElute, ScriptSeq v2, ScriptSeq Gold, RQ1 DNase and RQ1 Stop Solution) also showed the highest number of linked viral sequences (see Fig. 2).

In silico verification of viral sequences linked to laboratory components

Mapping of reads to reference genomes was conducted to identify genome coverage and to in silico validate results from the BLAST analysis. Both global and local alignments were performed. The coverage of reference genomes was reported using the global mapping results, whereas local mapping was a complement used to confirm local BLAST hits. Cross-library genome coverage of reference genomes above 80% was seen for 13% (63/493) of LCA viral sequences (see Supplementary material, table s5). Out of the 493 LCA viral sequences, 249 were linked to laboratory components

Table 2
Overview of the number of sequences analysed by BLAST

	Reads	Contigs
Sequences analysed by BLAST	2 953 972 594	1 381 107
Sequences identified by BLAST	790 424 528	574 477
Viral sequences identified by BLAST	91 863 018 (3.1%)	18 539 (1.3%)
Bleedover depleted viral sequences	91 654 946 (3.1%)	—
Bacterial sequences identified by BLAST	360 359 247 (12%)	411 889 (30%)
Other domain sequences identified by BLAST	338 202 263 (11%)	144 049 (10%)
Uncharacterized sequences	2 163 548 066 (73%)	806 628 (58%)

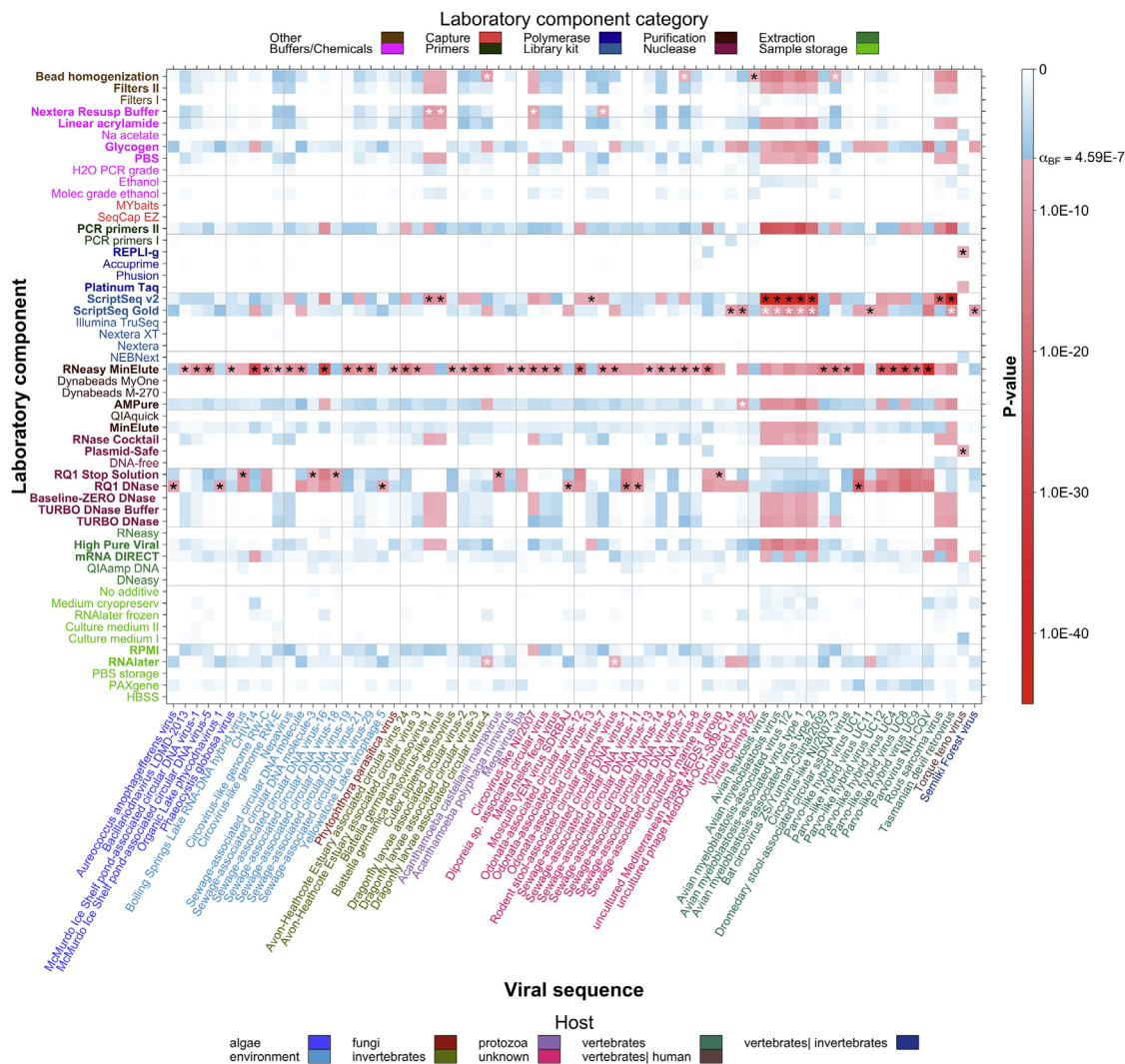


Fig. 1. p-values of association analysis between eukaryotic viral sequences and laboratory components at contig level. Including viral sequences linked to one or more laboratory components (e-value $<10^{-3}$). Significant associations illustrated in red and non-significant associations illustrated in blue. The strongest association(s) for each viral sequence is marked with a black star and white stars shows multiple-source associations. Laboratory components with minimum one linked viral sequence are marked in bold font.

based on global or local mapping of reads to reference genomes (see [Supplementary material, table s5](#)). Detailed investigation of viral clusters showed viral sequences composed of sequences proposed to originate from the identified virus (referred to as true viral sequences) and/or viral sequences assumed to originate from an unknown or non-viral source (referred to as artefact viral sequences). The artefact viral sequences were short and regionally repeated nucleotide sequences, generally of low complexity or showing homology to cloning vectors or human sequences.

Human host viral sequences

In total, 24 LCA viral sequences from viruses having humans as host were identified (see [Table 3](#) and [Supplementary material, table s4 and Results s1](#)). These viral sequences are particularly prone to erroneous conclusion when analysing human tissue samples. Low genome coverage (<25%) was identified in the majority of libraries (see [Supplementary material, fig. s5 and table s6](#)). A combination of sample-derived true viral sequences and laboratory-component-derived artefact sequences was identified for human mast-adenovirus C, human herpesvirus 1, human

herpesvirus 5, human immunodeficiency virus 1, human parvovirus B19 and torque teno virus. Among artefact sequences we identified homology to (i) various cloning or expression vectors (human immunodeficiency virus 1, human parvovirus B19 and Semliki Forest virus (see [Supplementary material, fig. s6 and Results s1](#))), (ii) human sequences (human papillomavirus type 6), and (iii) ribosomal RNA sequences (Simbu virus). Other artefact sequences did not show homology to specific types of sequences or were identified as short low complexity sequences.

Non-human vertebrate host viral sequences

We identified 60 viral sequences from viruses with a non-human vertebrate host among LCA viral sequences (see Fig. 1 and Supplementary material, fig. s2a and table s4). Among these, 29 were avian retroviruses (predominantly from the *Alpharetrovirus* genus), also including Tasmanian devil retrovirus (see Supplementary material, Results s2). The avian retroviral sequences were linked to ScriptSeq v2 and/or ScriptSeq Gold and were identified in high proportions (median above 60%) in libraries prepared using these methods (see Fig. 3). The cross-library genome

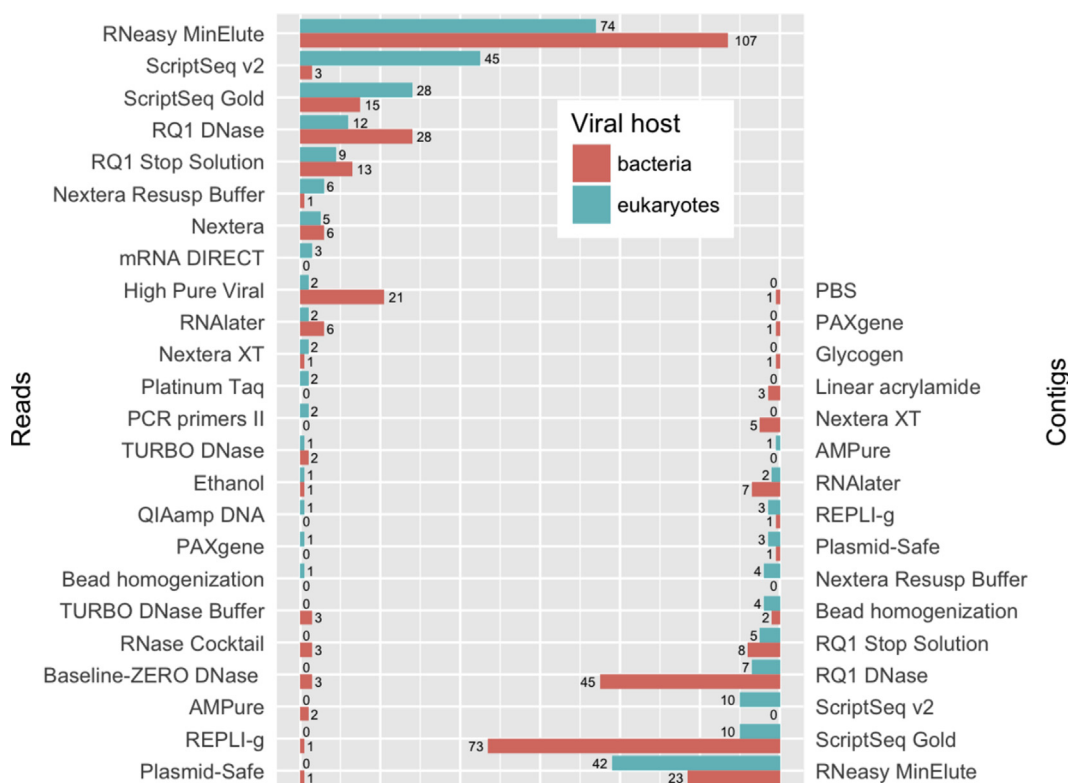


Fig. 2. Number of viral sequences linked to the various laboratory components. Counts comprise the number of viral sequences linked (showing the strongest association) to each laboratory component, including viral sequences linked to more than one laboratory component because of identical p-values, and multiple source viral sequences. Bars to the left and right show results from BLAST of reads and contigs, respectively.

coverage ranged from 4% to 100% (see [Supplementary material, table s5 and circos plots s1](#)) and 17 viral sequences showed coverage above 60% with dispersed alignments, therefore proposed to be true viral sequences originating from laboratory components. The remaining avian retroviral sequences are considered artefact viral sequences and true viral sequences present at low quantities. From the *Parvovirinae* subfamily 13 viral sequences were identified. These viral sequences showed cross-library genome coverage of 95%–100% and were all linked to the use of RNeasy MinElute, so were proposed to be laboratory-component-derived true viral sequences. Four viral sequences from the *Gammaretrovirus* genus linked to Nextera were identified. These showed regionally repeated alignments with relatively low cross-library genome coverage (7.8%–29%). The gammaretroviral sequences were, however, detected using several additional library preparation methods and we propose these sequences to be artefacts of unknown origin. In addition, 14 vertebrate viral sequences from eight different viral families were identified. Among these, Circovirus-like NI/2007-3 linked to RNeasy MinElute showed high cross-library genome coverage (96%), proposed to be a laboratory-component-derived true viral sequence. ASFV-like virus WU showed high cross-library genome coverage (80%) with dispersed alignments. It is considered a true viral sequence originating from laboratory components. The remaining vertebrate LCA viral sequences showed low coverage (<3%) with no or regionally repeated alignments, indicating artefact viral sequences originating from laboratory components.

Furthermore, a high number of LCA viral sequences from viruses with non-vertebrate hosts were identified (see [Fig. 1](#), and [Supplementary material, fig. s2a,b, and table s4](#)). Among these were 25 algal host viral sequences, containing 14 chlorella viruses

belonging to the *Chlorovirus* genus, including *Acanthocystis Turfacea Chlorella virus* (ATCV). All chlorella viral sequences were linked to RNeasy MinElute and showed no or low cross-library genome coverage (<5%) with dispersed and regionally repeated alignments (see [Supplementary material, table s5 and circos plots s1](#)). Chlorella viral sequences are proposed to be laboratory-component-derived artefacts and true viral sequences present at low quantities. The remaining algal viral sequences showed no or low cross-library genome coverage (<2%) with dispersed as well as regionally repeated alignments, indicating both artefact and true viral sequences originating from laboratory components. To investigate if the presence of algal viral sequences could be explained by the presence of algae, reads were globally mapped to six algal chloroplast genomes. The observed cross-library genome coverage was 6.4%–12% (see [Supplementary material, circos plots s2](#)). BLASTn of the mapped reads against the complete NCBI nucleotide database identified the same chloroplast genomes, thereby supporting the presence of algal sequences in our libraries. Moreover, 18 invertebrate, 14 environmental, 3 fungal, 8 plant, 13 protozoan, 28 unknown and 296 bacterial host viral sequences were identified (see [Supplementary material, Results s3](#)).

Non-template controls

Eight of the LCA viral sequences were not detected in any of the NTCs (see [Fig. 4a](#)). All eight were associated with the RNeasy MinElute kit. Among LCA viral sequences detected in the NTCs, the contaminating sequences were generally found in a higher proportion in NTC libraries than template-containing libraries. We can estimate the power to successfully detect the virus from the frequency of each specific virus in the NTCs. Taking avian

Table 3
Human viral sequences linked to laboratory components

Viral sequence	Blast reads/contigs association analysis		Mapping to reference genome		Evaluation of sequence origin and identity
	Linked laboratory component	p-value	N _{map}	Alignments global mapping (% coverage)	
Cyclovirus PK6197	RNeasy MinElute	2.5E-11	1	• 1 region 30 bp (1.7%)	LCD unknown artefact
Coxsackievirus B1	RQ1 Stop Solution	8.8E-09	0	—	LCD unknown artefact
Human mast-adenovirus C	QIAamp DNA	2.7E-10	151	• Dispersed (0.14%–85%) • 1 region <40 bp (<0.2%)	Sample-derived true viral sequences and LCD low complexity poly (A) artefact
Hepatitis E virus	ScriptSeq Gold	2.1E-12	0	—	LCD unknown artefact
Hepatitis C virus genotype 1	ScriptSeq v2	3.4E-19	537	• 4 regions <159 bp (<2%)	LCD low complexity poly (T) artefact
Hepatitis C virus subtype 1b	ScriptSeq v2	9.5E-10	533	• 4 regions <277 bp (<3%)	LCD low complexity poly (T) artefact
Human herpesvirus 1	PCR primers II	6.7E-11	81	• Dispersed (0.10%) • 3 regions <91 bp (<0.06%)	Sample-derived true viral sequences and LCD unknown artefact
Human herpesvirus 4	ScriptSeq v2	2.5E-25	12	• 2 regions <42 bp (0.03%)	LCD cloning vector artefact
Human herpesvirus 5	TURBO DNase	1.7E-21	43	• Dispersed (0.11%, 0.30%) • 1 region <208 bp (<0.1%)	Sample-derived true viral sequences and LCD cloning vector artefact
Human immunodeficiency virus 1	ScriptSeq v2	2.8E-33	113	• Dispersed (11%–67%) • 3 regions <205 bp (<3%)	Sample-derived true viral sequences and LCD cloning vector artefact
Human papillomavirus type 1a	Nextera XT	6.6E-10	16	• Dispersed (<4%)	True viral sequences of unknown origin
Human papillomavirus type 6	PAXgene	2.0E-07	8	• Dispersed (0.39%, 1.4%) • 3 regions <32 bp (<0.4%)	Human artefact
Human parvovirus B19	ScriptSeq v2	6.1E-12	41	• Dispersed (0.70%–100%) • 1 region <37 bp (<0.7%)	Sample-derived true viral sequences and LCD expression vector artefact
Human T-lymphotropic virus 1	ScriptSeq v2	2.1E-31	81	• 1 region <151 bp (<2%)	LCD unknown artefact
Influenza A virus ^a	Platinum Taq	2.8E-10	0	—	LCD unknown artefact
Influenza A virus ^b	ScriptSeq v2	1.2E-07	0	—	LCD unknown artefact
Influenza B virus ^c	Platinum Taq	1.2E-11	3	• Dispersed (1.5%) • 1 region <40 bp (<2%)	LCD unknown artefact
Lassa virus	ScriptSeq v2	7.8E-46	0	—	LCD unknown artefact
Merkel cell polyomavirus	Nextera XT	6.9E-17	34	• Dispersed (0.76%–76%)	LCD true viral sequences
Macaca mulatta polyomavirus 1	ScriptSeq Gold	1.2E-13	23	• Dispersed (11%) • 2 regions <93 bp (<2%)	Cross-mapping JC polyoma-virus and LCD expression and cloning vector artefact
Semliki Forest virus	ScriptSeq Gold	5.4E-10	0	—	LCD cloning vector artefact
Shamonda virus	RQ1 DNase	2.3E-11	0	—	LCD unknown artefact
Simbu virus	RQ1 DNase	3.2E-09	167	• 1 region <49 bp (0.40–0.70)	LCD ribosomal RNA artefact
Torque teno virus	Plasmid-Safe	1.9E-08	26	• Dispersed (3.2%–96%) • 1 region <351 bp (<10%)	Sample-derived true viral sequences and bleedover contamination

The table shows linked laboratory components and p-value of association analysis based on BLAST identification, number of libraries in which the viral sequence was identified by global mapping (N_{map}), distribution of alignments from the global mapping, evaluation of origin (sample-derived and/or laboratory-component-derived (LCD)) and type of sequence (true viral or artefact). For more detailed information see [Supplementary material, tables s4 and s5](#).

^a A/chicken/Karachi/NARC-100/2004(H7N3).

^b A/New York/55/2004(H3N2).

^c B/Thailand/CU-B2390/2010.

myeloblastosis virus and rodent stool-associated circular genome virus as examples; their respective detection frequencies in NTCs are 0.67 and 0.15. Assuming a binomial distribution, the probability of detecting these viral sequences if running three NTCs would be 0.96 and 0.39, respectively. To reach a probability of detection >0.95 for the rodent stool-associated circular genome virus, 19 NTCs would be necessary. [Fig. 4\(b\)](#) shows the number of NTCs needed for detection of a contaminating viral sequence in one or more NTCs with a 95% probability, illustrating the increasing number of NTCs necessary with decreasing detection rate.

Discussion

We have here provided a comprehensive list of 493 viral sequences, shown to be present in a variety of sample types and NTCs, significantly associated with the use of one or more laboratory components. Viral sequences showed stochastic appearance and were only detected in a subset of the libraries treated with the linked laboratory component, not always appearing in the NTCs. The hosts of linked viruses were taxonomically very diverse and included bacteria, protozoa, algae, plants, fungi, invertebrates and vertebrates.

To our knowledge, this is the first study using a systematic approach to identify a wide repertoire of contaminating viral

sequences and their origin. Several laboratory protocols and different laboratory components commonly used for sample preparation in virus discovery and surveillance with HTS were applied to the same samples, which facilitated the identification of laboratory-component-derived viral sequences. This is in contrast to other HTS studies where one laboratory practice has been applied to samples [28].

Viral sequence contamination in clinical samples is a known occurrence. Several viruses first linked to disease [19,21] have later been refuted as contamination [11,13,29,30]. In 2014, Yolken et al. linked ATCV-1 to altered cognitive function after its detection in the throat of healthy humans [20,31]. Subsequently, we refuted these findings and suggested that ATCV-1 corresponded to contamination arising from one or two laboratory components used concurrently during library preparation [17]. In 2014, a pipeline for identifying pathogens in HTS data from clinical samples was presented and applied to eight already published data sets originating from samples of various disease origins [32]. Noteworthy, 7.8% (76/974) of the viral findings in the study from 2014 were in this study linked to laboratory components.

Among the identified viral sequences, many lacked significant association to laboratory components. Many of these viruses are non-human and identification of these would not be expected in human cancer tissue samples. Some of the viral sequences have

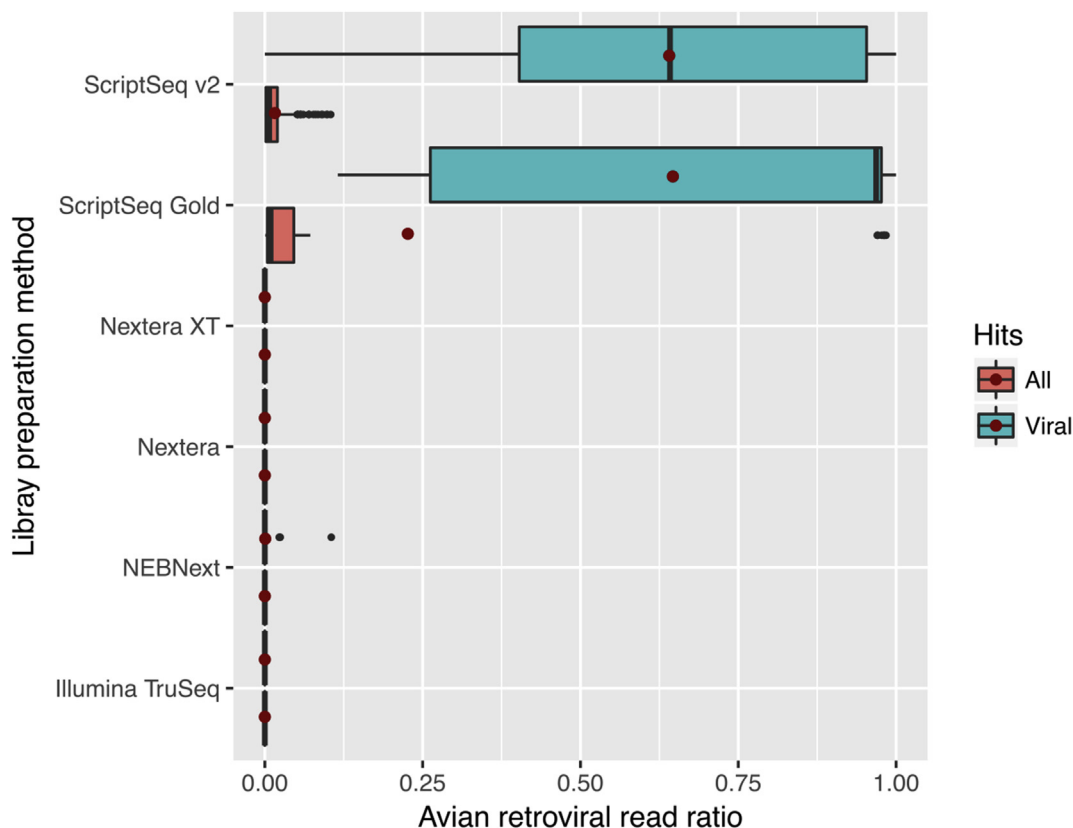


Fig. 3. Avian retroviral read ratios for the different library preparation methods. The black lines illustrate the median and the red dots illustrate the average ratio of avian retroviral reads for the different library preparation methods.

previously been suggested to be contamination, such as Pepino mosaic virus [33], Gallid herpes virus, Pandoravirus, Citrobacter phage [34] and Rotavirus [32,35,36].

More viral sequences significantly associated with laboratory components were identified among reads than contigs, a probable explanation being the low depth of coverage of genomes, making assembly of viral reads into contigs infrequent. Viral contigs were therefore only identified in a subset of libraries and more frequently when applying enrichment methods. This is a limitation of our study. As enrichment methods are expected to be better for detection of particularly low amounts of viral sequences, certain contaminating viral sequences may be statistically linked to enrichment-specific laboratory components because of their ability to detect them, though potentially having a different source. A notorious problem in HTS are reads occurring as a result of mixed sequencing clusters, bleedover. This phenomenon could explain the presence of specific viral sequences in NTCs (at ratios higher than the applied bleedover threshold). Our results indicate higher bleedover ratios in NTCs than in template-containing samples. Non-stringent e-values as low as 10^{-3} were applied in this study to reflect what is sometimes being applied in virus detection studies [37]. However, more stringent e-values are necessary (and are being applied) when using HTS to diagnose viral infection [38].

With regard to detection of novel contaminating viral sequences, we rely on some degree of sequence similarity in the BLAST identification. This is a limitation of the analysis and further effort could be put into identifying contaminating viral sequences in the unidentified sequences. A sequence recurrence-based clustering method has recently been published [15]. The strength of this approach is the independence of a sequence reference database for identification of correlation between nucleotide sequences and sample features.

Several viral sequences with significant association to several laboratory components were identified. Some of the laboratory components were used in parallel or almost in parallel, which resulted in identical or similar p-values, making it hard to be certain of the origin. The factual origin(s) of viral sequences could be verified by setting up a designed experiment with different combinations of laboratory components and using virus-specific PCR. This was, however, beyond the scope of this study.

Independent of the bioinformatic method used for the initial identification of viral sequences, we find it of outmost importance to evaluate genome coverage, read depth and the distribution of alignments across the identified viral genome. A high coverage and/or dispersed sequences across the reference genome indicate that the viral sequences are derived from the virus in question rather than representing an artefact. Short regionally repeated viral sequences in multiple samples indicate artefact viral sequences and should always raise suspicion. Re-blasting of regionally repeated LCA viral sequences showed additional best hits to cloning and expression vectors, for example, Semliki Forest virus, used as a vector for vaccine development, for gene therapy and for production of recombinant proteins [39–41]. Others showed no additional best hit and the artefact sequence could therefore not be identified. The wide use of viral cloning and expression vectors could be an overlooked problem in virus discovery, leading to false positives.

Concerning the suitability of NTCs, our main conclusion is that several negative controls should be included in order to detect sporadic contaminants, despite the costs of sequencing (see [Supplementary material, Discussion s1](#)). Furthermore, we strongly recommend that viral sequences from viruses with non-human hosts should be handled with caution when identified in HTS data and that researchers carefully consider the possibility of contamination. It should be noted that viral sequences that we

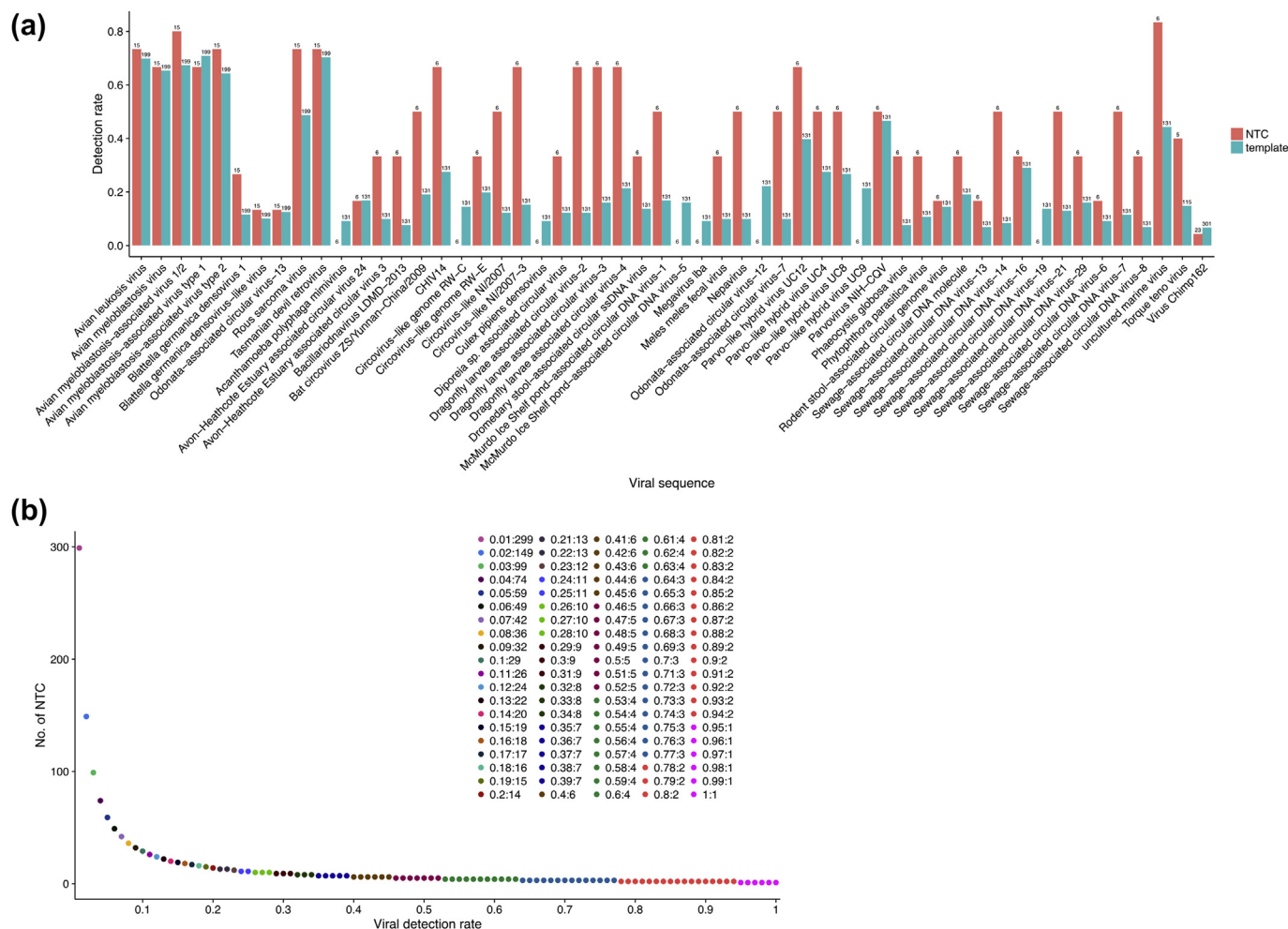


Fig. 4. Non-template controls. (a) Detection rates of specific laboratory-component-associated viral sequences in non-template control (NTC) libraries and template-containing libraries. The number of NTC libraries and template-containing libraries is shown above the bars in the figure. (b) Number of NTCs necessary to reach a detection probability of minimum 0.95 for different viral detection rates.

propose as contamination in human tissue samples, could have a natural origin in samples from another species or environmental location.

Transparency declaration

The authors declare no conflicts of interest.

Parts of the data have previously been presented in Geneva October 2018 at the International Conference on Clinical Metagenomics (ICCMg).

Funding

This work was supported by Innovation Fund Denmark grant No 019-2011-2 (The Genome Denmark platform) and Danish National Research Foundation grant No DNRF94. The funders had no role in study design, data collection, analysis and interpretation, decision to submit the work for publication, or preparation of the manuscript.

Authors' contributions

MA and KRK wrote the manuscript. SM and JMGI made major revisions of the manuscript. Laboratory experiments were designed

by SM, HF, LV, KRK and RHJ and performed by SM, KRK, HF, LV, IBN, SRR, RHJ, ARI, DEAP and PVSO. Study design was done by AJH, JMGI, KRK, MA, TM, SM, LPN, OL, SB, TSP and EW. MA, JFN, JMGI and TAH designed the bioinformatic pipeline. Initial bioinformatic analysis (pre-processing, assembly and BLAST) was conducted by JFN, JMGI and MLMH. MA and KRK performed the concluding analysis (clustering, data mining, statistical analysis, visualization). Mapping analysis and creation of Circos plots was done by JARH.

Acknowledgement

We thank BGI Europe for sequencing of the samples.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.cmi.2019.04.028>.

References

- Reuter JA, Spacek DV, Snyder MP. High-throughput sequencing technologies. *Mol Cell* 2015;58:586–97.
- Oulas A, Pavlou C, Polymenakou P, Pavlopoulos GA, Papanikolaou N, Kotoulas G, et al. Metagenomics: tools and insights for analyzing next-

- generation sequencing data derived from biodiversity studies. *Bioinform Biol Insights* 2015;9:75–88.
- [3] Briesse T, Paweska JT, McMullan LK, Hutchison SK, Street C, Palacios G, et al. Genetic detection and characterization of Lujo virus, a new hemorrhagic fever-associated arenavirus from southern Africa. *PLoS Pathogens* 2009;5:e1000455.
 - [4] Yozwiak NL, Skewes-Cox P, Gordon A, Saborio S, Kuan G, Balmaseda A, et al. Human enterovirus 109: a novel interspecies recombinant enterovirus isolated from a case of acute pediatric respiratory illness in Nicaragua. *J Virol* 2010;84:9047–58.
 - [5] Feng H, Shuda M, Chang Y, Moore PS. Clonal integration of a polyomavirus in human Merkel cell carcinoma. *Science* 2008;319:1096–100.
 - [6] Cholleti H, Hayer J, Abilio AP, Mulandane FC, Verner-Carlsson J, Falk KI, et al. Discovery of novel viruses in mosquitoes from the Zambezi Valley of Mozambique. *PLoS One* 2016;11:e0162751.
 - [7] Hansen TA, Fridholm H, Frøsløv TG, Kjartansdóttir KR, Willerslev E, Nielsen LP, et al. New type of papillomavirus and novel circular single stranded DNA virus discovered in urban *Rattus norvegicus* using circular DNA enrichment and metagenomics. *PLoS One* 2015;10:e0141952.
 - [8] Ng TFF, Manire C, Borrowman K, Langer T, Ehrhart L, Breitbart M. Discovery of a novel single-stranded DNA virus from a sea turtle fibropapilloma by using viral metagenomics. *J Virol* 2009;83:2500–9.
 - [9] Munang'andu HM, Mugimba KK, Byarugaba DK, Mutoloki S, Evensen Ø. Current advances on virus discovery and diagnostic role of viral metagenomics in aquatic organisms. *Front Microbiol* 2017;8:406.
 - [10] Chiu CY. Viral pathogen discovery. *Curr Opin Microbiol* 2013;16:468–78.
 - [11] Smuts H, Kew M, Khan A, Korsman S. Novel hybrid parvovirus-like virus, NIH-CQV/PHV, contaminants in silica column-based nucleic acid extraction kits. *J Virol* 2014;88: 1398–8.
 - [12] Lusk RW. Diverse and widespread contamination evident in the unmapped depths of high throughput sequencing data. *PLoS One* 2014;9:e110808.
 - [13] Naccache SN, Greninger AL, Lee D, Coffey LL, Phan T, Rein-Weston A, et al. The perils of pathogen discovery: origin of a novel parvovirus-like hybrid genome traced to nucleic acid extraction spin columns. *J Virol* 2013;87:11966–77.
 - [14] Lysholm F, Wetterbom A, Lindau C, Darban H, Bjerkner A, Fahlander K, et al. Characterization of the viral microbiome in patients with severe lower respiratory tract infections, using metagenomic sequencing. *PLoS One* 2012;7:e30875.
 - [15] Friis-Nielsen J, Kjartansdóttir KR, Møllerup S, Asplund M, Mourier T, Jensen RH, et al. Identification of known and novel recurrent viral sequences in data from multiple patients and multiple cancers. *Viruses* 2016;8:53.
 - [16] Laurence M, Hatzis C, Brash DE. Common contaminants in next-generation sequencing that hinder discovery of low-abundance microbes. *PLoS One* 2014;9:e97876–8.
 - [17] Kjartansdóttir KR, Friis-Nielsen J, Asplund M, Møllerup S, Mourier T, Jensen RH, et al. Traces of ATCV-1 associated with laboratory component contamination. *Proc Natl Acad Sci USA* 2015;112:E925–6.
 - [18] Lo S-C, Pripuzova N, Li B, Komaroff AL, Hung G-C, Wang R, et al. Detection of MLV-related virus gene sequences in blood of patients with chronic fatigue syndrome and healthy blood donors. *Proc Natl Acad Sci USA* 2010;107: 15874–9.
 - [19] Lombardi VC, Ruscetti FW, Gupta Das J, Pfost MA, Hagen KS, Peterson DL, et al. Detection of an infectious retrovirus, XMRV, in blood cells of patients with chronic fatigue syndrome. *Science* 2009;326:585–9.
 - [20] Yolken RH, Jones-Brando L, Dunigan DD, Kannan G, Dickerson F, Severance E, et al. Chlorovirus ATCV-1 is part of the human oropharyngeal virome and is associated with changes in cognitive functions in humans and mice. *Proc Natl Acad Sci USA* 2014;111:16106–11.
 - [21] Xu B, Zhi N, Hu G, Wan Z, Zheng X, Liu X, et al. Hybrid DNA virus in Chinese patients with seronegative hepatitis discovered by deep sequencing. *Proc Natl Acad Sci USA* 2013;110:10264–9.
 - [22] Schlager R, Choe DJ, Brown KR, Thaker HM, Singh IR. XMRV is present in malignant prostatic epithelium and is associated with prostate cancer, especially high-grade tumors. *Proc Natl Acad Sci USA* 2009;106:16351–6.
 - [23] Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucl Acids Res* 1997;25:3389–402.
 - [24] Kircher M, Sawyer S, Meyer M. Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucl Acids Res* 2012;40: e3–e3.
 - [25] R Core Team. R: a language and environment for statistical computing [Internet]. Vienna, Austria. Available at: <https://www.R-project.org/>.
 - [26] Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Meth* 2012;9:357–9.
 - [27] Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, et al. Circos: an information aesthetic for comparative genomics. *Genome Res* 2009;19: 1639–45.
 - [28] Moustafa A, Xie C, Kirkness E, Biggs W, Wong E, Turpaz Y, et al. The blood DNA virome in 8,000 humans. *PLoS Pathogens* 2017;13:e1006292.
 - [29] Erlwein O, Robinson MJ, Dustan S, Weber J, Kaye S, McClure MO. DNA extraction columns contaminated with murine sequences. *PLoS One* 2011;6: e23484.
 - [30] Paprotka T, Delviks-Frankenberry KA, Cingöz O, Martinez A, Kung H-J, Tepper CG, et al. Recombinant origin of the retrovirus XMRV. *Science* 2011;333:97–101.
 - [31] Yolken RH, Jones-Brando L, Dunigan DD, Kannan G, Dickerson F, Severance E, et al. Reply to Kjartansdóttir et al.: Chlorovirus ATCV-1 findings not explained by contamination. *Proc Natl Acad Sci USA* 2015;112: E927–7.
 - [32] Naccache SN, Federman S, Veeraraghavan N, Zaharia M, Lee D, Samayoa E, et al. A cloud-compatible bioinformatics pipeline for ultrarapid pathogen identification from next-generation sequencing of clinical samples. *Genome Res* 2014;24:1180–92.
 - [33] Bukowska-Oško I, Perlejewski K, Nakamura S, Motooka D, Stokowy T, Kosińska J, et al. Sensitivity of next-generation sequencing metagenomic analysis for detection of RNA and DNA viruses in cerebrospinal fluid: the confounding effect of background contamination. *Adv Exp Med Biol* 2016. epub ahead of print.
 - [34] Hjelmso MH, Hellmér M, Fernandez-Cassi X, Timoneda N, Lukjancenko O, Seidel M, et al. Evaluation of methods for the concentration and extraction of viruses from sewage in the context of metagenomic sequencing. *PLoS One* 2017;12:e0170199.
 - [35] Grard G, Fair JN, Lee D, Slikas E, Steffen I, Muyembe J-J, et al. A novel rhabdovirus associated with acute hemorrhagic fever in Central Africa. *PLoS Pathogens* 2012;8:e1002924.
 - [36] Grard G, Fair JN, Lee D, Slikas E, Steffen I, Muyembe J-J, et al. Correction: a novel rhabdovirus associated with acute hemorrhagic fever in Central Africa. *PLoS Pathogens* 2016;12:e1005503.
 - [37] Prachayaprecha S, Schapendonk CME, Koopmans MP, Osterhaus ADME, Schürch AC, Pas SD, et al. Exploring the potential of next-generation sequencing in detection of respiratory viruses. *J Clin Microbiol* 2014;52: 3722–30.
 - [38] Théze J, Li T, Plessis du L, Bouquet J, Kraemer MUG, Somasekar S, et al. Genomic epidemiology reconstructs the introduction and spread of Zika virus in Central America and Mexico. *Cell Host Microbe* 2018;23:855–7.
 - [39] Atkins GJ, Fleeton MN, Sheahan BJ. Therapeutic and prophylactic applications of alphavirus vectors. *Expert Rev Mol Med* 2008;10:e33.
 - [40] Lundström K. Alphavirus vectors as tools in neuroscience and gene therapy. *Virus Res* 2016;216:16–25.
 - [41] DiCiommo DP, Bremner R. Rapid, high level protein production using DNA-based Semliki Forest virus vectors. *J Biol Chem* 1998;273:18060–6.